

Jacobian 神经网络算法

甄景贤 (King-Yin Yan)

General.Intelligence@Gmail.com

Abstract.

经典的神经网络 Back Prop 学习算法，它是一个 error-driven 算法，但在很多人工智能的实际应用中，不存在唯一的「理想答案」，而是根据正或负的奖励 (reward) 学习。当答案正确时，奖励 > 0 , error = 0; 当答案不正确时，奖励 < 0 ，但 error 仍是不知道的（因为不知道理想答案）。简言之，就是不能用 error-driven 学习。

所以我想出了一个 reward-driven 的学习法：假设神经网络将 $\mathbf{x}_0 \mapsto \mathbf{y}_0$ ，它通常也会将 \mathbf{x}_0 的邻域 map 到 \mathbf{y}_0 的邻域。如果我们想「加强」这个映射，可以将「更大的 \mathbf{x}_0 的邻域」映射到「接近 \mathbf{y}_0 的邻域」。

这种算法对人工智能应该很重要，暂时我还想不出有什么其他办法，可以做到 [深度] 神经网络的 reward-driven 学习。

将这思想更准确化，可以将 feed-forward 神经网络的构造看成是这样的：

$$\mathbf{y} = \mathbf{F}(\mathbf{x}) \quad (1)$$

$$\mathbf{y} = \bigcirc^L \mathbf{W} \dots \bigcirc^\ell \mathbf{W} \dots \bigcirc^1 \mathbf{W} \mathbf{x} \quad (2)$$

其中 \mathbf{W} 代表每一层 (layer) ℓ 的矩阵。

\mathbf{F} 的反方向是：

$$\mathbf{x} = \mathbf{F}^{-1}(\mathbf{y}) \quad (3)$$

$$\mathbf{x} = \overset{1}{\rhd} \ominus \dots \overset{\ell}{\rhd} \ominus \dots \overset{L}{\rhd} \ominus \mathbf{y} \quad (4)$$

注意： $\overset{1}{\rhd} = \mathbf{W}^{-1}$ ， $\ominus = \bigcirc^{-1}$ ，形状不同。

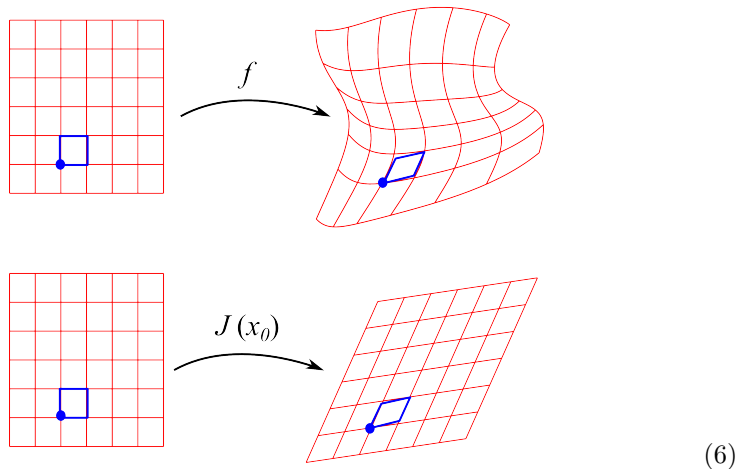
假设在 \mathbf{x} 空间有体积元 U ，经过 \mathbf{F} 变换成 \mathbf{y} 空间的体积元 V ，那么：

$$U = |J| \cdot V \quad (5)$$

$J = \left[\frac{\partial \mathbf{F}(\mathbf{x})}{\partial \mathbf{x}} \right]$ 叫 Jacobian 矩阵。

在我们的情况下， $|J| = \left| \frac{\partial \mathbf{F}^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right|$ 在 \mathbf{y}_0 的值，代表「单位体积元由 $\mathbf{y}_0 \mapsto \mathbf{x}_0$ 的变化率」。（下面会看到， \mathbf{F} 和 \mathbf{F}^{-1} 的正/反方向不太重要，因为基本上不影响计算复杂度。）

Jacobian $J(\mathbf{x}_0)$ is a local linearization of the function f near the point \mathbf{x}_0



每次得到**正奖励**，我们会令 Jacobian $|J|$ 增加一点：

$$|J| := \det \left[\frac{\partial \mathbf{F}^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right]_{n \times n} \quad (7)$$

下标表示那是一个 $n \times n$ 矩阵。

其实 Jacobian 矩阵的意义就是：

$$J = \left[\frac{\partial \text{输出}}{\partial \text{输入}} \right] \quad (8)$$

神经网络的输入和输出都是 $\dim n$ ，所以 Jacobian 很自然是 $n \times n$ 矩阵。

用**梯度下降法**，我们需要计算这些梯度： $\left[\frac{\partial |J|}{\partial \mathbf{W}} \right]$ ，总数是网络中的 weights 的个数 $= \sum m_\ell$ 。

要用到 determinant 的微分公式：

$$\frac{d}{dt} |A(t)| = \text{tr}(\text{adj}(A) \cdot \frac{dA(t)}{dt}) \quad (9)$$

$$\text{adj}(A) := |A| \cdot A^{-1} \quad (10)$$

换句话说，对於每个权重 $w := \bar{W}_{ij}^\ell$ ，我们要计算：

$$\frac{\partial}{\partial w} |J| = \text{tr}(|J| \cdot J^{-1} \cdot \left[\frac{\partial J}{\partial w} \right]) \quad (11)$$

注意： $|J|$ 和 J^{-1} 是 \mathbf{y}_0 的函数，只需在大 loop 外一次过计算。

问题是，计算 $\left[\frac{\partial J}{\partial w} \right]_{n \times n}$ 的时候：

$$\frac{\partial J}{\partial w} = \frac{\partial}{\partial w} \frac{\partial \mathbf{F}^{-1}}{\partial \mathbf{y}} = \frac{\partial}{\partial w} \frac{\partial}{\partial \mathbf{y}} \stackrel{1}{\geq} \bigcirc \dots \stackrel{\ell}{\geq} \bigcirc \dots \stackrel{L}{\geq} \bigcirc \mathbf{y} \quad (12)$$

这牵涉到用 w 对 W^{-1} 的分量微分，可以想像就算计了出来也会是极复杂的。解决办法是，索性「本末倒置」，用 \geq 来定义神经网络，然后在 forward propagation 时才用 $W = \geq^{-1}$ 计算。

J 的分量写出来是：

$$J_{ij} = \frac{\partial \mathbf{F}_i^{-1}}{\partial y_j} = \frac{\partial}{\partial y_j} \left[\stackrel{1}{\geq} \bigcirc \dots \stackrel{\ell}{\geq} \bigcirc \dots \stackrel{L}{\geq} \bigcirc \mathbf{y} \right]_i =: \nabla_{ij}^1 \quad (13)$$

$$\begin{cases} \nabla_{ij}^1 &:= \sum_{k_1} \left[\stackrel{1}{\geq}_{ik_1} \bigcirc' (y_{k_1}^2) \nabla_{ij}^2 \right] \\ \nabla_{ij}^\ell &:= \sum_{k_\ell} \left[\stackrel{\ell}{\geq}_{k_{\ell-1}k_\ell} \bigcirc' (y_{k_\ell}^{\ell+1}) \nabla_{ij}^{\ell+1} \right] \\ \nabla_{ij}^L &:= \stackrel{L}{\geq}_{k_{L-1}j} \bigcirc' (y_j) \end{cases} \quad (14)$$

这情况完全类似於经典 Back Prop，以上只是 chain rule 的应用， ∇^ℓ 将每层用 chain rule 分拆开来，所以 ∇ 又叫“local gradient”。上式就是整个网络的**反向传递**，其中每个 weight 出现 exactly 一次。

但工作还未完，我们要计算 $\frac{\partial J_{ij}}{\partial \geq} = \dot{\nabla}_{ij}^1$ 。（定义 $\geq := \stackrel{\ell}{\geq}_{gh}$ ， $k_0 := i$ ， $k_L := j$ ）

注意： $\mathbf{x} = \mathbf{F}^{-1}(\mathbf{y})$ ，所以 \mathbf{y} 是自变量， \geq 不影响 \mathbf{y} ，所以 $\frac{\partial \mathbf{y}}{\partial \geq} \equiv 0$ 。

\geq 必会是 $\stackrel{\ell}{\geq}$ 的其中一元，但如果 $\geq \notin \stackrel{\ell}{\geq}$ ，以下的项微分后都会变成 0：

$$\begin{cases} \dot{\nabla}_{ij}^1 = \sum_{k_1} \left[\stackrel{1}{\geq}_{ik_1} \bigcirc' (y_{k_1}^2) \dot{\nabla}_{ij}^2 \right] \\ \dot{\nabla}_{ij}^\ell = \sum_{k_\ell} \left[\stackrel{\ell}{\geq}_{k_{\ell-1}k_\ell} \bigcirc' (y_{k_\ell}^{\ell+1}) \dot{\nabla}_{ij}^{\ell+1} \right] \\ \dot{\nabla}_{ij}^L = \stackrel{L}{\geq}_{k_{L-1}j} \bigcirc' (y_j) \equiv 0 \end{cases} \quad (15)$$

所以实际上只剩下一项：

$$\frac{\partial J_{ij}}{\partial \succ} = \sum_{k_1} \left[\overset{1}{\succ}_{k_0 k_1} \odot'(y_{k_1}^2) \dots \sum_{k_{\ell-2}} \left[\overset{\ell-2}{\succ}_{k_{\ell-3} k_{\ell-2}} \odot'(y_{k_{\ell-2}}^{\ell-1}) \dots \right. \right. \quad (16)$$

$$\left. \left. \begin{cases} \dots \overset{\ell-1}{\succ}_{k_{\ell-2}, g} \odot'(y_g^\ell) \odot'(y_h^{\ell+1}) \nabla_{ij}^{\ell+1} \end{cases} \right] \right]$$

$$\left. \begin{cases} \dots \overset{\ell-1}{\succ}_{k_{\ell-2}, g} \odot'(y_g^\ell) \odot'(y_h^{\ell+1}) \end{cases} \right] \quad \text{if } \succ \in \text{last layer}$$

上式的意思是：每层 layer 重复一块 $[\sum \succ \odot]$ ，直到遇到 $\succ = \overset{\ell}{\succ}_{gh}$ ，则用结尾形式取代之。

和经典 Back Prop 不同的是，上式只是 $n \times n$ 矩阵中的一个元素，从复杂度而论，每个 weight 的 ∇ 计算，增加了起码 n^2 倍的复杂度（虽然其计算上可以共用一些结果）。记住 $n = \dim$ 状态空间。

可以这样理解：每个 weight 的调教，需要计算这个 weight 对 Jacobian 的影响，而那 Jacobian 是**整个网络**的特性。关键似乎就在於每个 weight 对 Jacobian 的影响。

现在回看更高层次的这个式子：

$$\frac{\partial}{\partial w} |J| = \text{tr}(|J| \cdot J^{-1} \cdot \left[\frac{\partial J}{\partial w} \right]) \quad (17)$$

$$= |J| \cdot \text{tr} \left(\left[\frac{\partial y}{\partial x} \right] \cdot \left[\frac{\partial}{\partial w} \frac{\partial x}{\partial y} \right] \right) \quad (18)$$

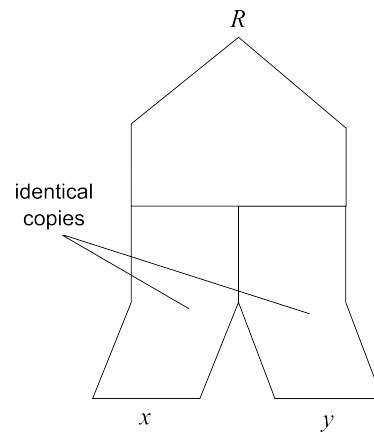
$$= |J| \cdot \sum_{ij} \left(\frac{\partial y}{\partial x} \right)_{ij} \left(\frac{\partial}{\partial w} \frac{\partial x}{\partial y} \right)_{ij} \quad (19)$$

上式中最重要（最慢）的是那 $(i, j) \in n \times n$ 求和。裡面的第一个因子是 Jacobian J ，第二个因子是我们刚计算了的 $\nabla_w J^{-1}$ 。

Back Prop 的 ∇ 形式上是 $\frac{\partial \text{输出}}{\partial w}$ ，我们的 ∇ 形式是 $\left[\frac{\partial}{\partial w} \frac{\partial \text{输入}}{\partial \text{输出}} \right]_{n \times n}$ 。

其实我们只需要计算 $\nabla_w |J|$ 的**大约方向**。暂时我在代码中的做法是：忽略式 (19) 中较小的项，那就不需做足 n^2 个乘积。

或者可不可以将 $|J|(\succ)$ 看成是一个 weight \succ 的函数，然后用它的 Taylor series expansion 来近似？



(20)

或者，照旧是 feedforward network，但 somehow 它的学习是基于某 reward function。问题是这 reward function 从何而来？ R 可以是一个全部 W 的函数，是由我们任意定义的。

1 Jacobian learning algorithm

The classical Back Prop algorithm is **error-driven**, but in many AI problems the “correct” answers are not given, instead the feedback is provided via **rewards**. When an answer is correct, the reward $R > 0$, error $\mathcal{E} = 0$; when the answer is incorrect, $R < 0$, but \mathcal{E} is still unknown (because we don’t know the correct answer). In other words, error-driven learning is inapplicable.

So I thought of a reward-driven learning method: assume the neural network maps $\mathbf{x}_0 \mapsto \mathbf{y}_0$, usually it also maps the neighborhood of \mathbf{x}_0 to the neighborhood of \mathbf{y}_0 . If we wish to “strengthen” this pair of mapping, we can make a **bigger** neighborhood of \mathbf{x}_0 map to the same neighborhood close to \mathbf{y}_0 . We will make this precise.

This type of learning algorithm should be very useful to AI, and currently the author is not aware of other alternatives for training [deep] neural networks via rewards.

A feed-forward neural network can be constructed this way:

$$\mathbf{y} = \mathbf{F}(\mathbf{x}) \quad (21)$$

$$\mathbf{y} = \bigcirc \overset{\ell}{W} \dots \bigcirc \overset{\ell}{W} \dots \bigcirc \overset{1}{W} \mathbf{x} \quad (22)$$

where W represents the matrix of **weights** on each layer ℓ .

The **inverse** of F is:

$$\mathbf{x} = \mathbf{F}^{-1}(\mathbf{y}) \quad (23)$$

$$\mathbf{x} = \overset{1}{\oslash} \dots \overset{\ell}{\oslash} \dots \overset{L}{\oslash} \mathbf{y} \quad (24)$$

Note: $\oslash = W^{-1}$, $\bigcirc = \bigcirc^{-1}$, the shape is different.

Assume that in the space of \mathbf{x} there is a volume element U , which transforms via \mathbf{F} to a volume element V in the space of \mathbf{y} , then:

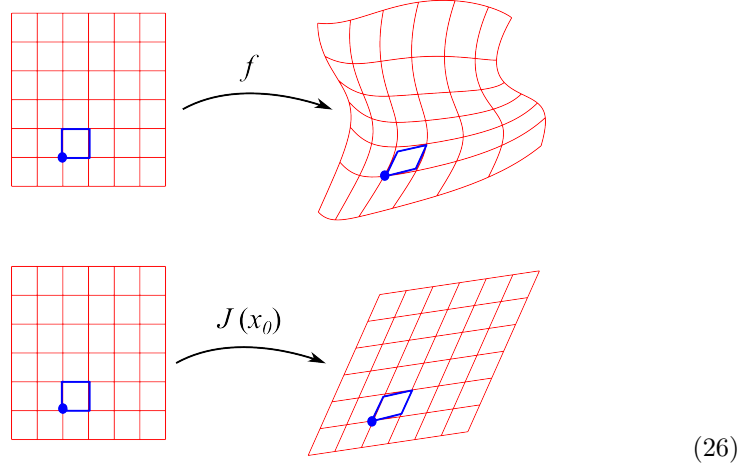
$$U = |J| \cdot V \quad (25)$$

$J = \left[\frac{\partial \mathbf{F}(\mathbf{x})}{\partial \mathbf{x}} \right]$ is called the **Jacobian** matrix.

In our case, the value of $|J| = \left| \frac{\partial \mathbf{F}^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right|$ at \mathbf{y}_0 represents “the change in volume from $\mathbf{y}_0 \mapsto \mathbf{x}_0$ ”. (Below, we will see that the direction of \mathbf{F} or \mathbf{F}^{-1} is not too

important, as either way the computational complexity is essentially the same.)

Jacobian $J(x_0)$ is a local linearization of the function f near the point x_0



Every time we get a **positive reward**, we can let the Jacobian $|J|$ increase slightly:

$$|J| := \det \left[\frac{\partial F^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right]_{n \times n} \quad (27)$$

The subscript indicates that it is a $n \times n$ matrix.

The meaning of the Jacobian matrix is:

$$J = \left[\frac{\partial \text{input}}{\partial \text{output}} \right] \quad (28)$$

The neural network's input and output are both of $\dim n$, so the Jacobian is naturally an $n \times n$ matrix.

To use **gradient descent**, we need to calculate these gradients: $\left[\frac{\partial |J|}{\partial \mathbf{W}} \right]$, their total number is the number of weights in the network $= \sum \ell \#(\mathbf{W})$.

We need the formula for the derivative of the determinant:

$$\frac{d}{dt} |A(t)| = \text{tr}(\text{adj}(A) \cdot \frac{dA(t)}{dt}) \quad (29)$$

$$\text{adj}(A) := |A| \cdot A^{-1} \quad (30)$$

In other words, for each weight $w := \overset{\ell}{W}_{ij}$, we need to calculate:

$$\frac{\partial}{\partial w} |J| = \text{tr}(|J| \cdot J^{-1} \cdot \left[\frac{\partial J}{\partial w} \right]) \quad (31)$$

Note: $|J|$ and J^{-1} are functions of \mathbf{y}_0 , we only need to calculate them once outside the big loop.

Now a problem arises in the calculation of $\left[\frac{\partial J}{\partial w} \right]_{n \times n}$:

$$\frac{\partial J}{\partial w} = \frac{\partial}{\partial w} \frac{\partial \mathbf{F}^{-1}}{\partial \mathbf{y}} = \frac{\partial}{\partial w} \frac{\partial}{\partial \mathbf{y}} \stackrel{1}{\rhd} \ominus \dots \stackrel{\ell}{\rhd} \ominus \dots \stackrel{L}{\rhd} \ominus \mathbf{y} \quad (32)$$

This requires us to differentiation W^{-1} w.r.t. w ; we can imagine the result would be very complicated. So we use a trick, by “reversing” the network, we use \rhd to define the network weights, and then use $W = \rhd^{-1}$ during forward propagation.

The components of J are:

$$J_{ij} = \frac{\partial \mathbf{F}_i^{-1}}{\partial y_j} = \frac{\partial}{\partial y_j} \left[\stackrel{1}{\rhd} \ominus \dots \stackrel{\ell}{\rhd} \ominus \dots \stackrel{L}{\rhd} \ominus \mathbf{y} \right]_i =: \nabla_{ij}^1 \quad (33)$$

$$\begin{cases} \nabla_{ij}^1 &:= \sum_{k_1} \left[\stackrel{1}{\rhd}_{ik_1} \ominus'(y_j^2) \nabla_{ij}^2 \right] \\ \nabla_{ij}^\ell &:= \sum_{k_\ell} \left[\stackrel{\ell}{\rhd}_{k_{\ell-1}k_\ell} \ominus'(y_j^{\ell+1}) \nabla_{ij}^{\ell+1} \right] \\ \nabla_{ij}^L &:= \stackrel{L}{\rhd}_{k_{L-1}j} \ominus'(y_j) \end{cases} \quad (34)$$

This situation is exactly analogous to the classical Back Prop algorithm; The above is just the application of the **chain rule**, with ∇^ℓ written separately for each layer, therefore ∇ is called the “local gradient”. The above formula amounts to propagating the entire network one time, where every weight appears **exactly once**.

But our work is not finished yet; We need to calculate $\frac{\partial J_{ij}}{\partial \rhd} =: \dot{\nabla}_{ij}^1$.

(Let’s define $\rhd := \stackrel{\ell}{\rhd}_{gh}$, $k_0 := i$, $k_L := j$)

\rhd would be an element of \rhd , but if $\rhd \notin \rhd$, all the terms below would vanish:

$$\begin{cases} \dot{\nabla}_{ij}^1 = \sum_{k_1} \left[\stackrel{1}{\rhd}_{ik_1} \ominus'(y_j^2) \dot{\nabla}_{ij}^2 \right] \\ \dot{\nabla}_{ij}^\ell = \sum_{k_\ell} \left[\stackrel{\ell}{\rhd}_{k_{\ell-1}k_\ell} \ominus'(y_j^{\ell+1}) \dot{\nabla}_{ij}^{\ell+1} \right] \\ \dot{\nabla}_{ij}^L = \stackrel{L}{\rhd}_{k_{L-1}j} \ominus'(y_j) \equiv 0 \end{cases} \quad (35)$$

So what is left over is just this term:

$$\frac{\partial J_{ij}}{\partial \mathfrak{z}} = \sum_{k_1} \left[\mathfrak{z}_{k_0 k_1}^1 \odot'(y_j^2) \dots \sum_{k_\ell} \left[\mathfrak{z}_{k_{\ell-1} k_\ell}^\ell \odot'(y_j^{\ell+1}) \dots \right. \right. \quad (36)$$

$$\left. \left. \begin{cases} \dots \odot'(y_j^{\ell+1}) \nabla_{ij}^{\ell+1} \\ \dots \odot'(y_j^{\ell+1}) \end{cases} \right] \right] \quad \text{if } \mathfrak{z} \in \text{last layer}$$

The above formula means: For each layer we repeat a block of $[\sum \mathfrak{z} \odot]$, until we encounter $\mathfrak{z} = \mathfrak{z}_{gh}^\ell$, then we replace with the terminal form.

In contrast to classical Back Prop, the above formula gives us only one element in an $n \times n$ matrix; From the complexity point of view, calculating the ∇ for each weight is at least n^2 times as costly as classical Back Prop (even though we may re-use some intermediate computation results). Recall that $n = \dim \boxed{\text{state space}}$.

We can understand it thusly: For each weight we try to calculate its influence towards the Jacobian, but the Jacobian is a **global** property of the network. The key seems to lie in how each weight **influences** the Jacobian.

Now let's look back at this higher-level formula:

$$\frac{\partial}{\partial w} |J| = \text{tr}(|J| \cdot J^{-1} \cdot \left[\frac{\partial J}{\partial w} \right]) \quad (37)$$

$$= |J| \cdot \text{tr} \left(\left[\frac{\partial y}{\partial x} \right] \cdot \left[\frac{\partial}{\partial w} \frac{\partial x}{\partial y} \right] \right) \quad (38)$$

$$= |J| \cdot \sum_{ij} \left(\frac{\partial y}{\partial x} \right)_{ij} \left(\frac{\partial}{\partial w} \frac{\partial x}{\partial y} \right)_{ij} \quad (39)$$

The most critical (slowest) part is the $(i, j) \in n \times n$ summation. The first factor inside \sum is the Jacobian J , the second factor is the $\nabla_w J^{-1}$ that we just calculated.

Back Prop's ∇ has the form $\frac{\partial \boxed{\text{output}}}{\partial \boxed{\text{weights}}}$

whereas our ∇ has the form $\left[\frac{\partial}{\partial \boxed{\text{weights}}} \frac{\partial \boxed{\text{input}}}{\partial \boxed{\text{output}}} \right]_{n \times n}$.

In fact we just need to calculate the **approximate** direction and size of $\nabla_w |J|$. Currently in our code we use this trick: ignore the smaller terms in (39), so we don't need to do all of n^2 products.

Or perhaps we can regard $|J|(\mathfrak{z})$ as a function of the weight \mathfrak{z} , and then use its Taylor series expansion to approximate?